

Scalability of Graph Neural Networks

Prof. O-Joun Lee

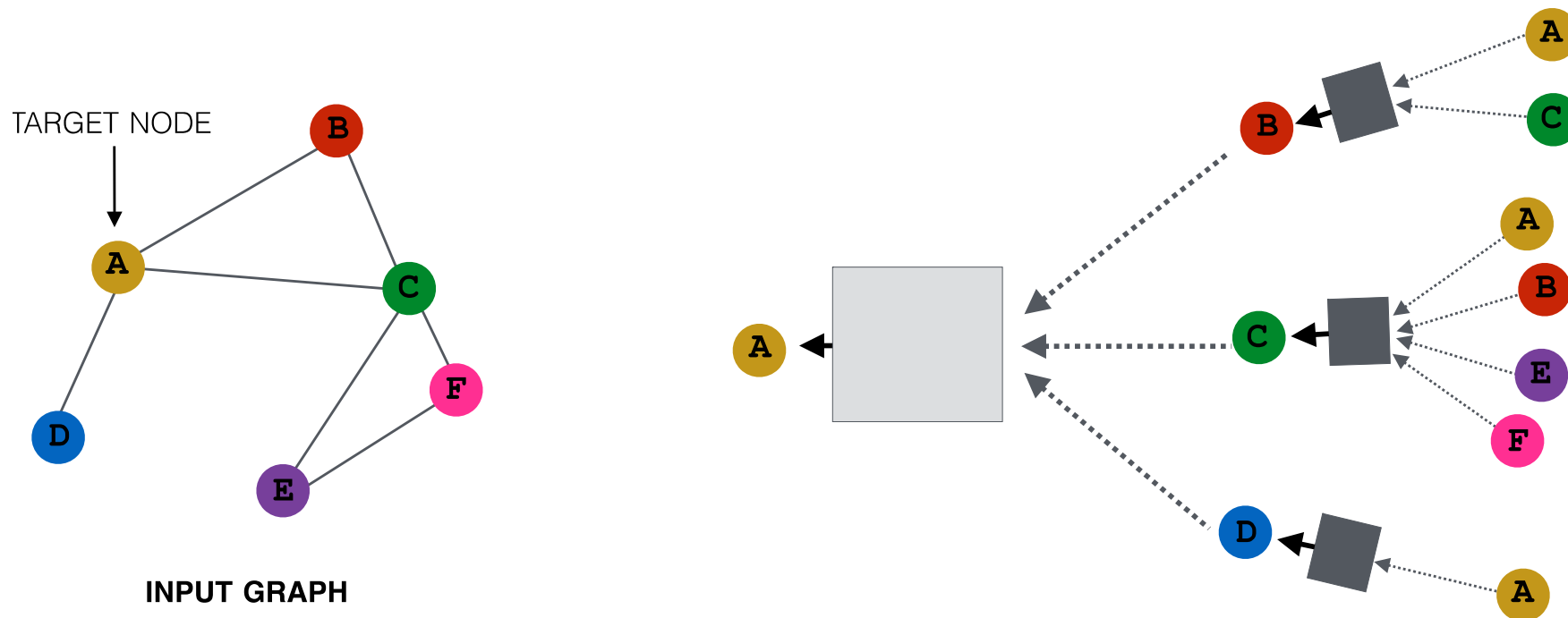
Dept. of Artificial Intelligence,
The Catholic University of Korea
ojlee@catholic.ac.kr

Contents

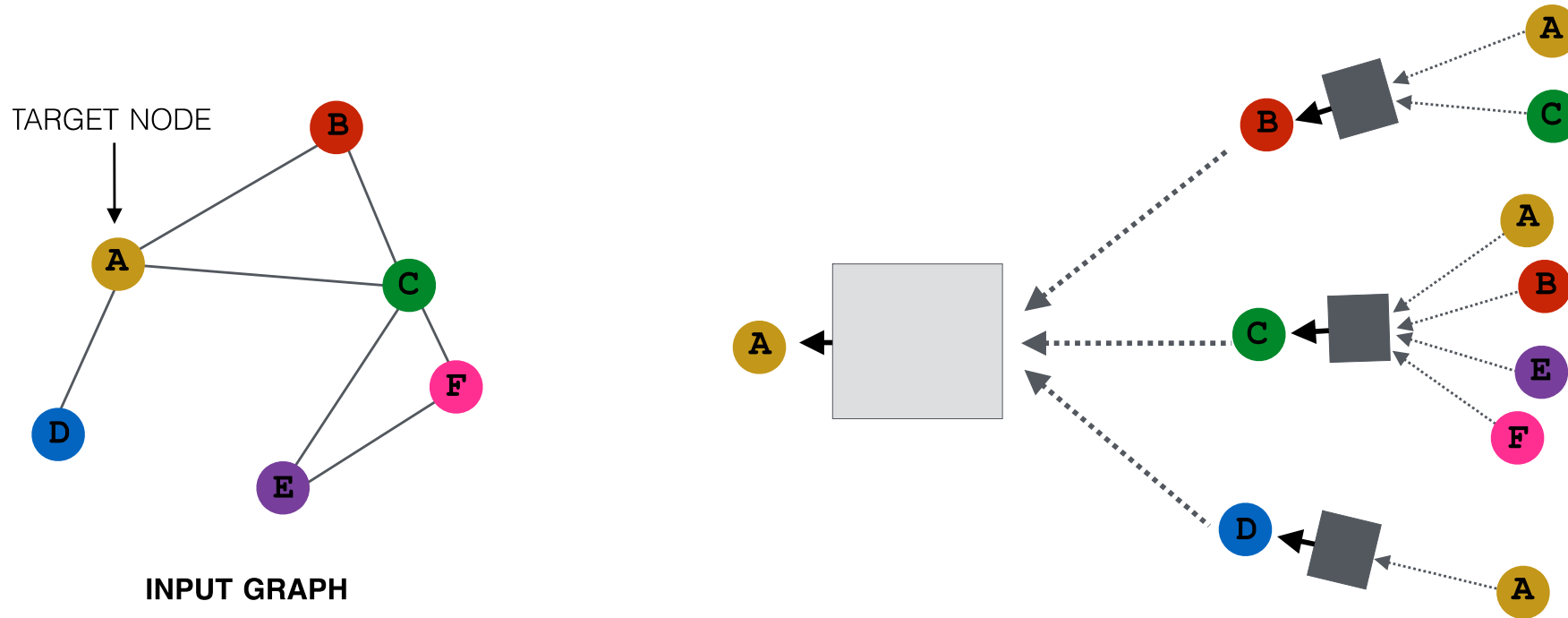


- Issues towards the large-scale GNNs
- Node-wise sampling with GraphSage
- Graph-wise sampling with ClusterGCN
- GraphSAINT

- **Key idea:** Generate node embeddings based on local neighborhoods.

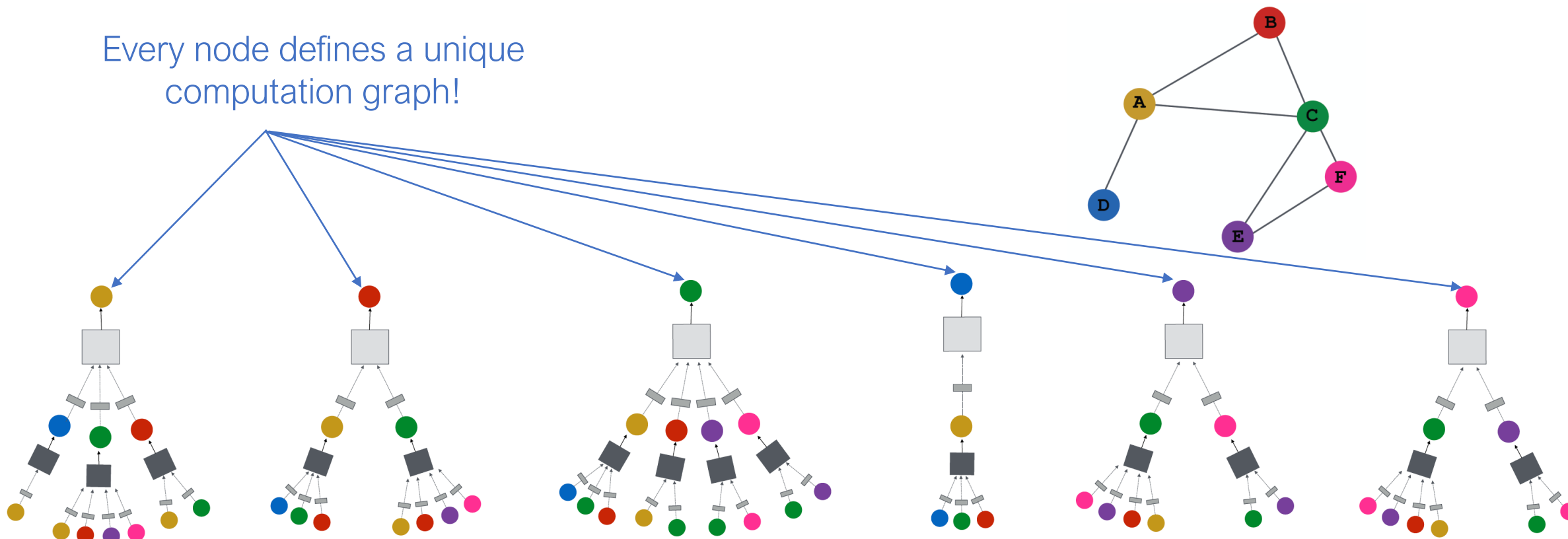


- **Intuition:** Nodes aggregate information from their neighbors using neural networks

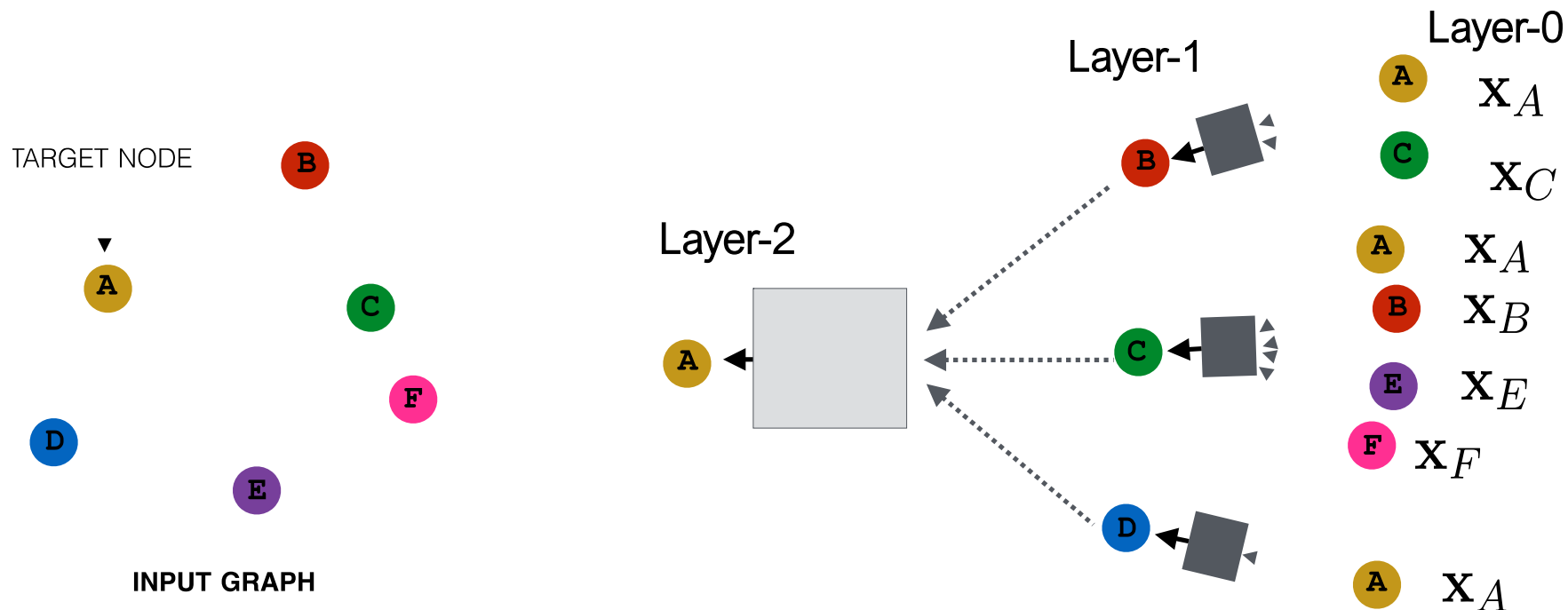


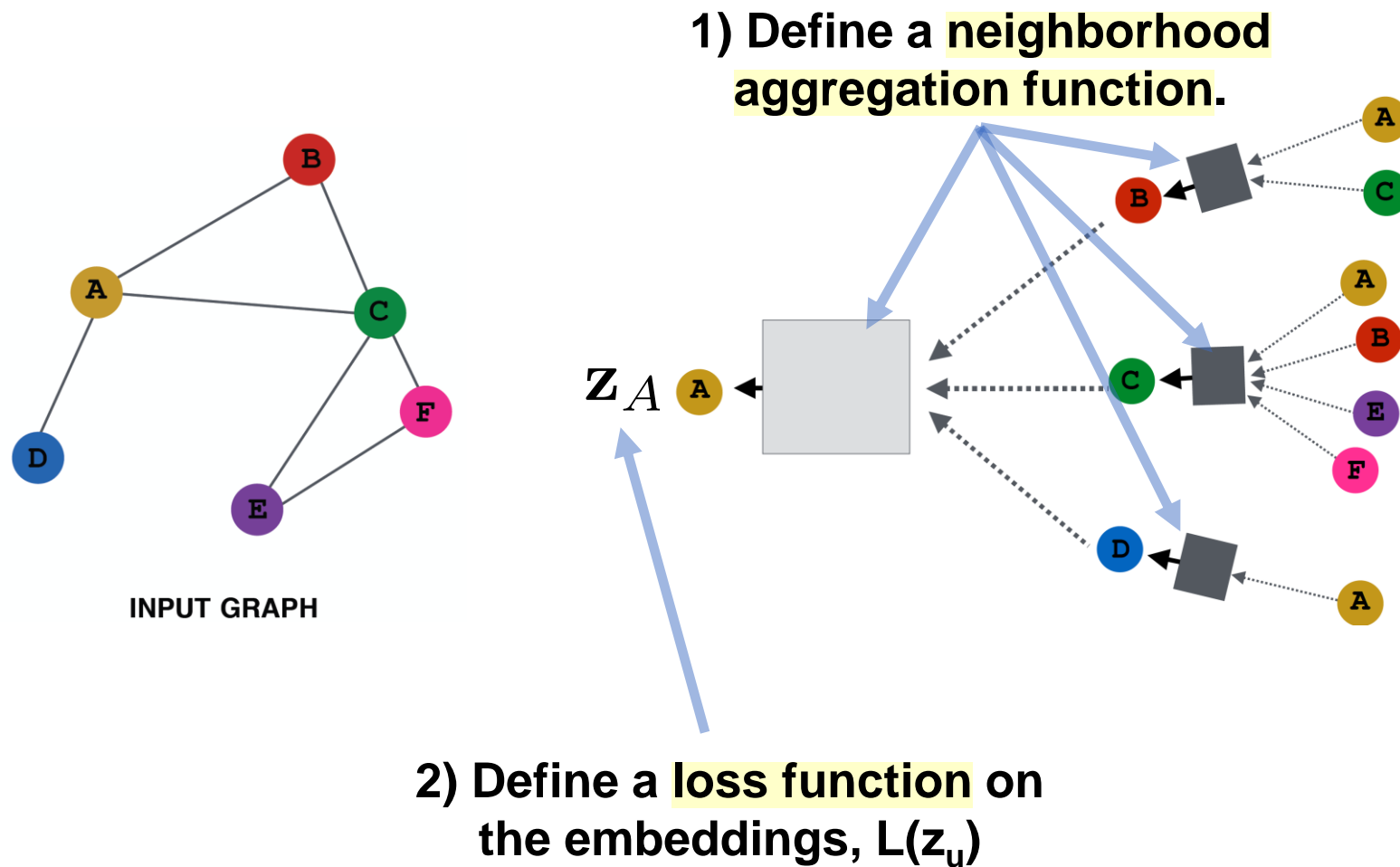
- **Intuition:** Network neighborhood defines a computation graph

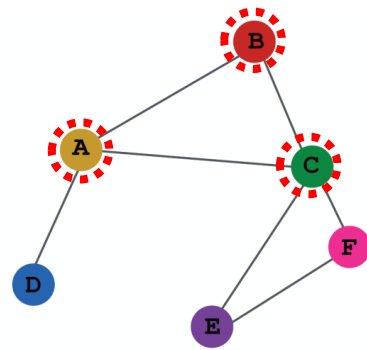
Every node defines a unique computation graph!



- Nodes have **embeddings** at each layer.
- Model can be arbitrary depth.
- “layer-0” embedding of node A is its input feature, i.e. x_A

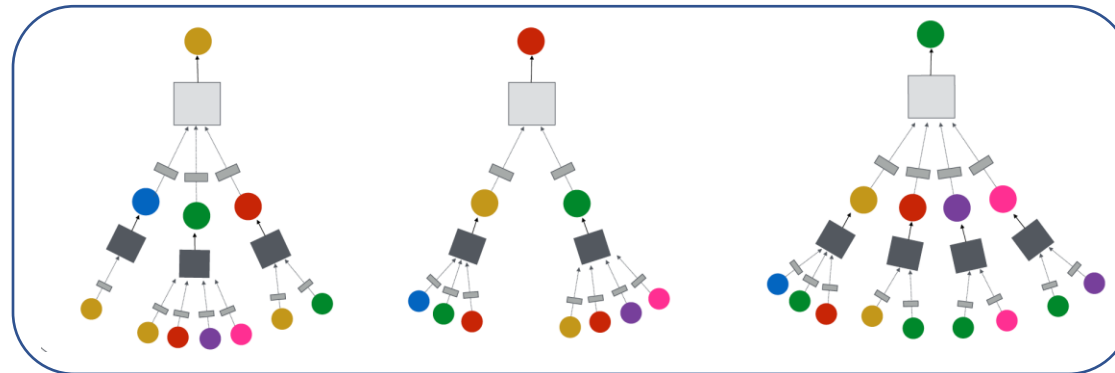






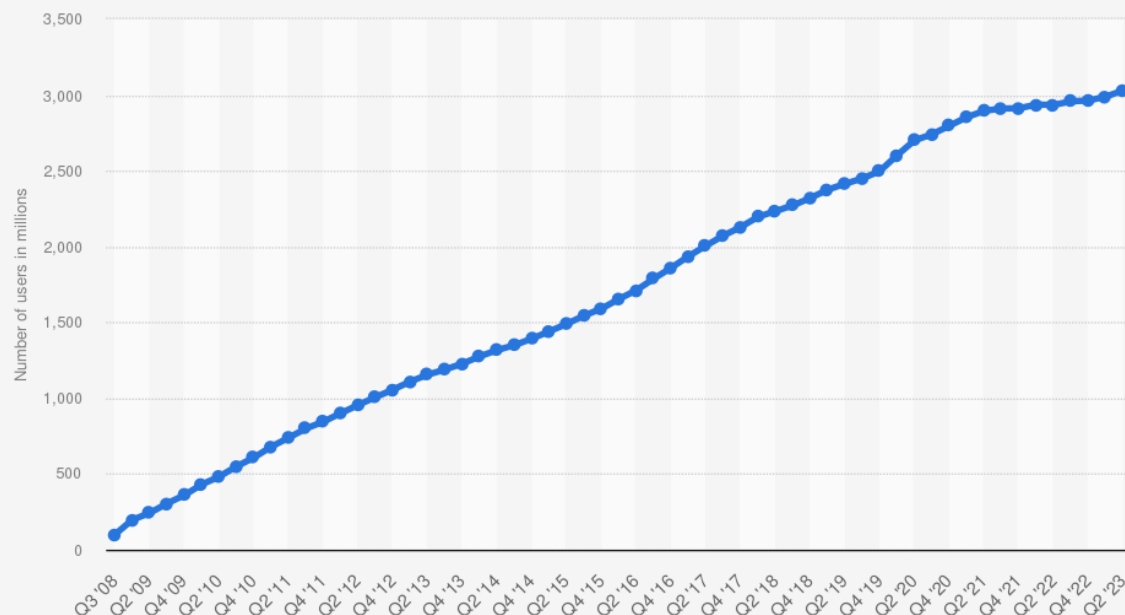
INPUT GRAPH

3) Train on a set of nodes, i.e., a batch of compute graphs



Large scale

Number of monthly active Facebook users worldwide as of 2nd quarter 2023 (in millions)



Source
Meta Platforms
© Statista 2023

Additional Information:
Worldwide; Meta Platforms; Q3 2008 to Q2 2023

Large number

ZINC Substances Catalogs Tranches Biological More

ZINC15

Welcome to ZINC, a free database of commercially-available compounds for virtual screening. ZINC contains over 230 million purchasable compounds in ready-to-dock, 3D formats. ZINC also contains over 750 million purchasable compounds you can search for analogs in under a minute.

Getting Started

- [Getting Started](#)
- [What's New](#)
- [About ZINC 15 Resources](#)
- [Current Status / In Progress](#)
- [Why are ZINC results "estimates"?](#)

Ask Questions

You can use ZINC for **general** questions such as

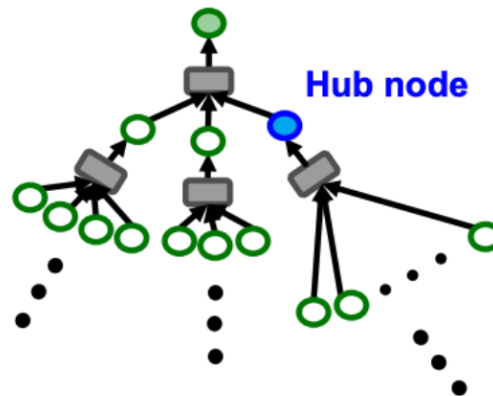
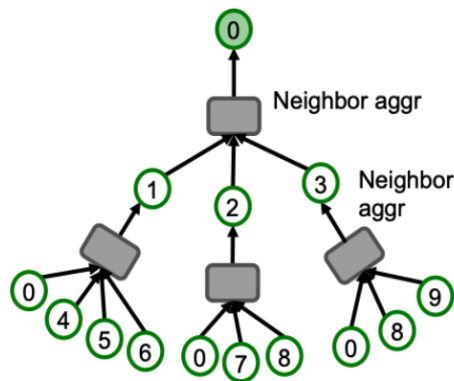
- [How many substances in current clinical trials have PALL](#)
- [How many natural products have names in ZINC and ar SMILES, names and calculated logP](#)
- [How many endogenous human metabolites are there? \(](#)

➤ **Computationally Expensive:**

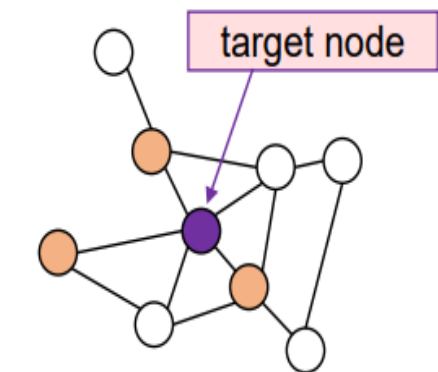
- We need to generate the complete K-hop neighborhood computational graph and then need to aggregate plenty of information from its surroundings.
- As we go deeper into the neighborhood computation graph becomes exponentially large.
 - problem while fitting these computational graphs inside GPU memory.

➤ **The curse of Hub nodes or Celebrity nodes:**

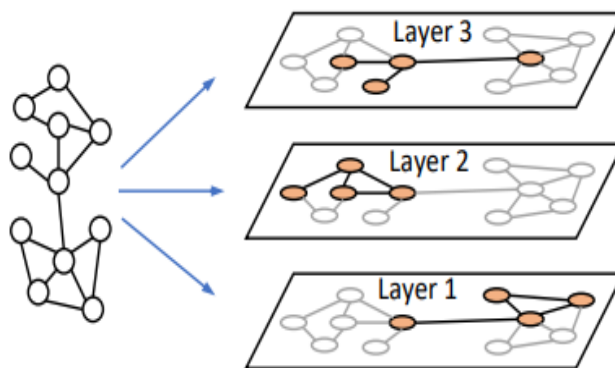
- Hub nodes are those nodes which are very high degree nodes in the graph



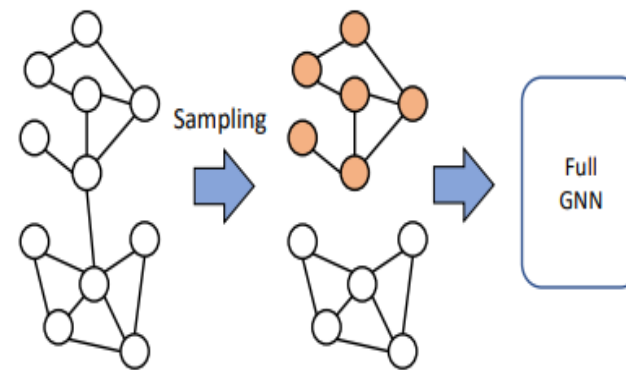
- **Why the original GNN fails on large graph?**
 - Large memory requirement.
 - Inefficient gradient update.
- **Three paradigms toward large-scale GNN:**
 - Node-wise sampling
 - Layer-wise sampling
 - Graph-wise sampling



Node-wise Sampling

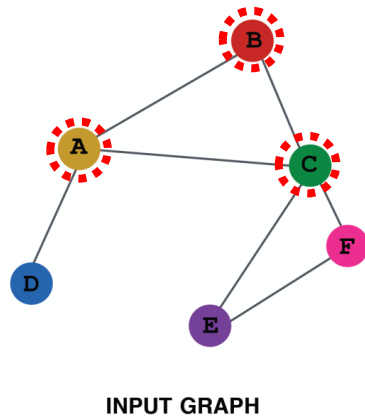


Layer-wise Sampling

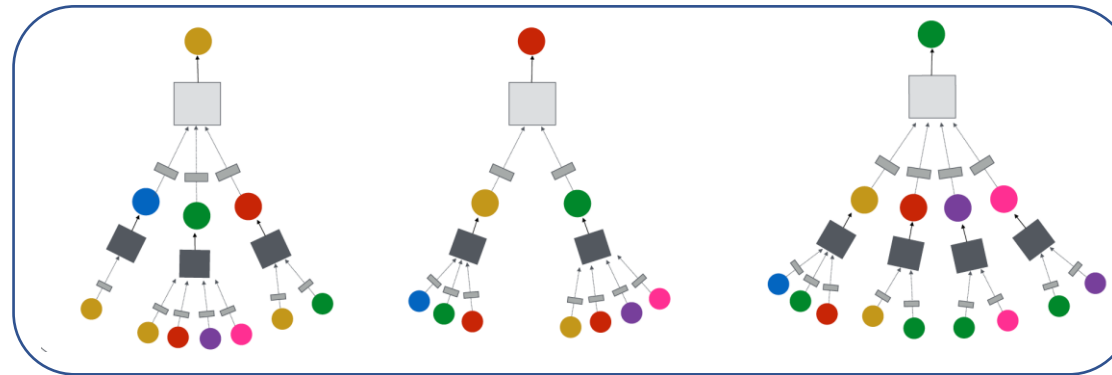


Graph-wise Sampling

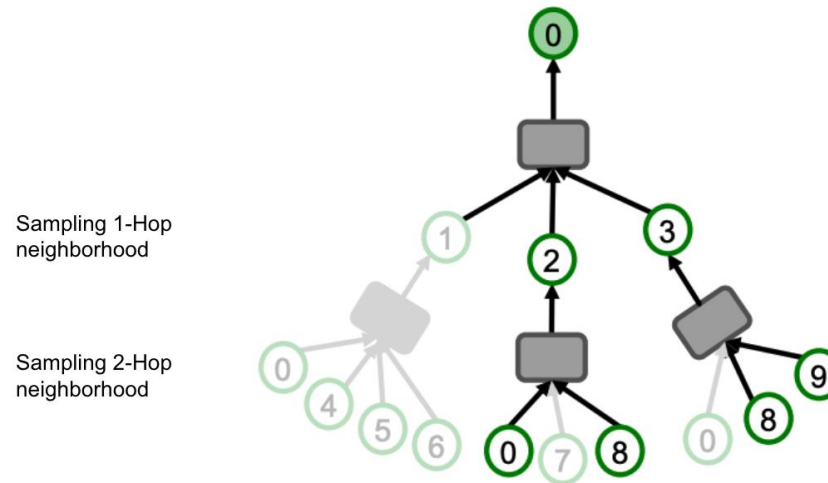
- How to design efficient sampling algorithm?
- How to guarantee the sampling quality?



Train on a set of nodes, i.e., a batch of compute graphs

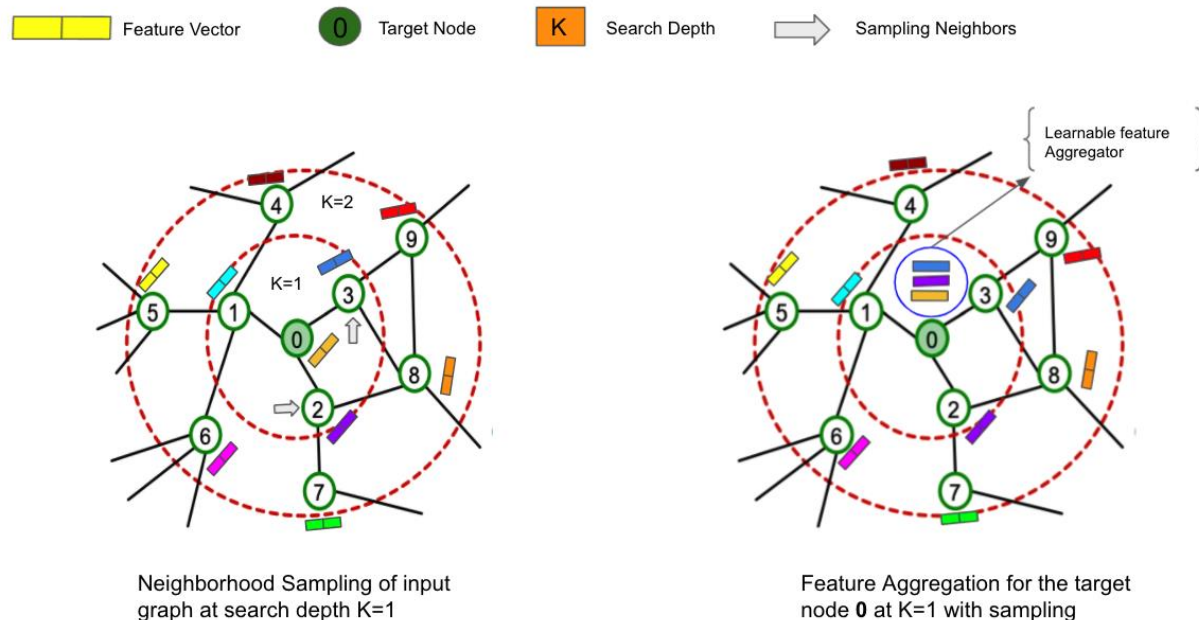


- So far, we have aggregated the neighbor messages by taking their (weighted) average, can we do better?
- **The idea:** not take the entire K-hop neighborhood of a target node but **select few nodes at random from the K-hop neighborhood** in order to generate computational graph.
- This process is known as neighborhood sampling which provides the GraphSage algorithm its unique ability of scaling up to billion of nodes in the graph.



- GraphSage is an inductive version of GCNs which implies that it does not require the whole graph structure during learning and it can generalize well to the unseen nodes.
- We don't need to learn the embeddings for each node.
- Learning an aggregation function (MEAN, POOLING, or LSTM) which when given an information (or features) from the local neighborhood of a node then it knows how to aggregate those features (learning takes place via stochastic gradient descent)

GraphSage



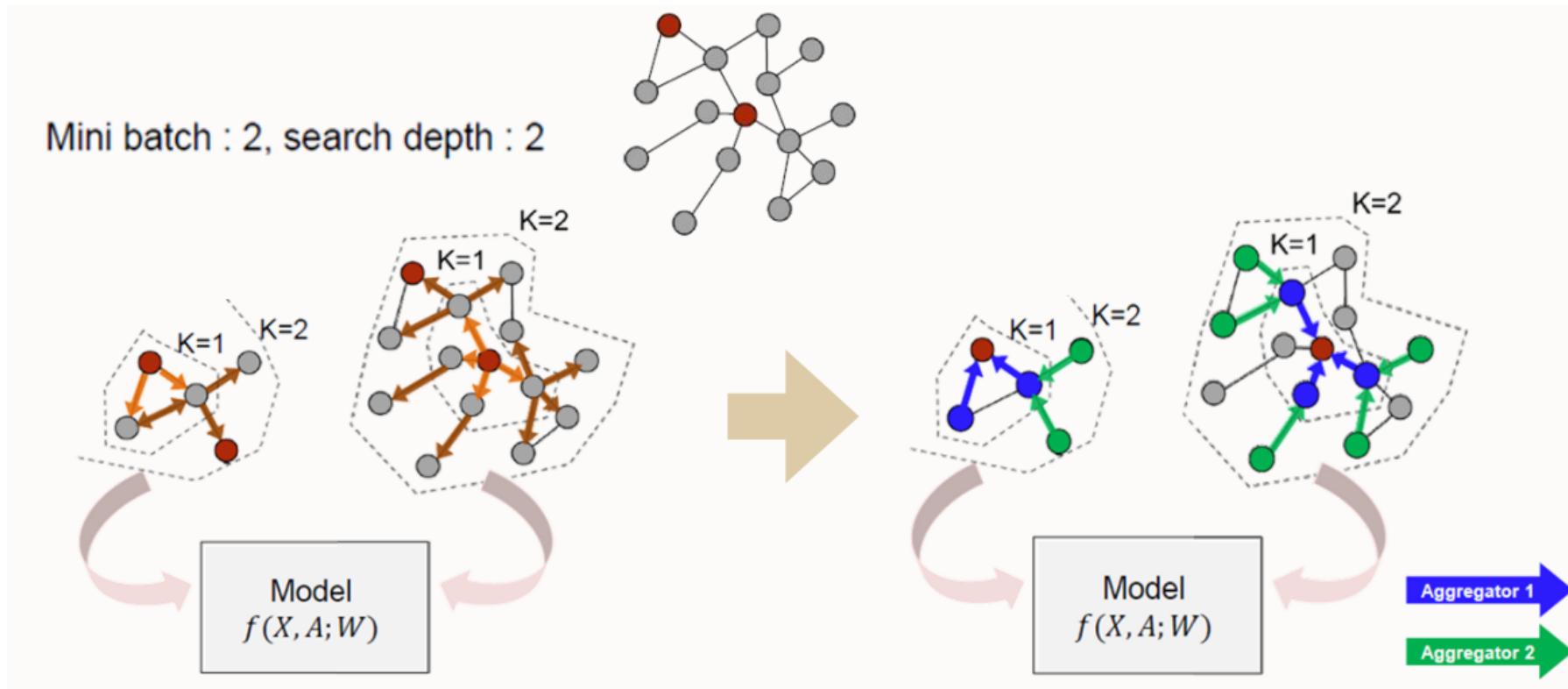
- GraphSAGE (SAmple and aggreGatE)
 - Instead of training individual embeddings for each node, generates embeddings by sampling features from neighborhoods

→ Mini batch

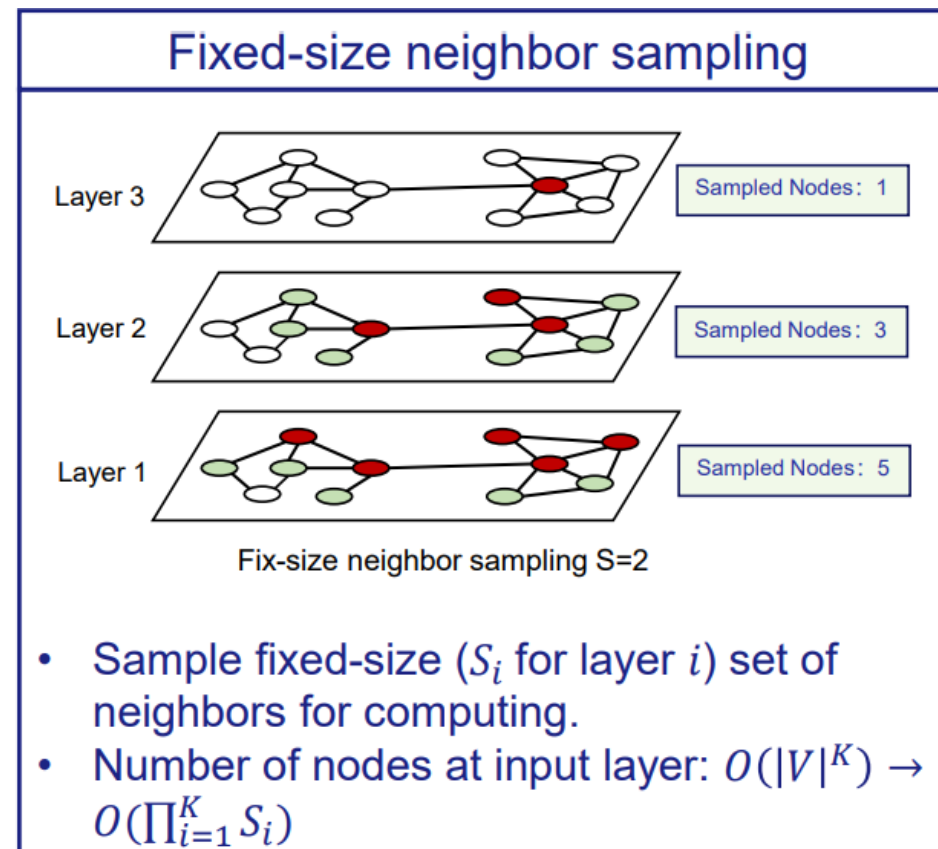
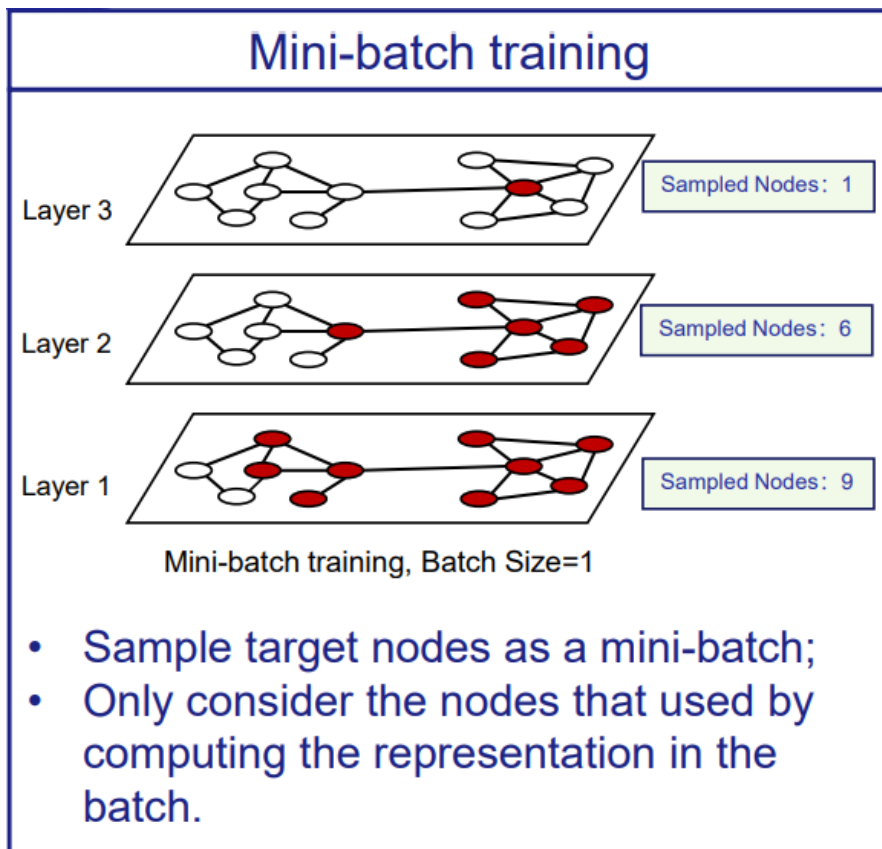
- Train a set of aggregator functions that learn to aggregate feature information a node's local neighborhood

→ Aggregating

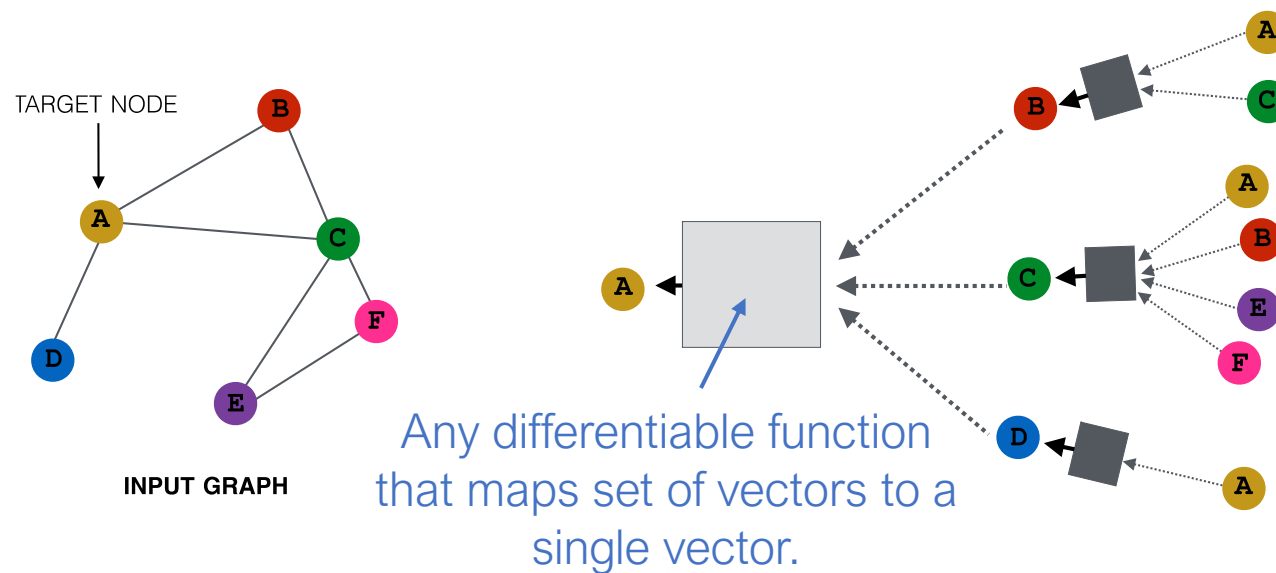
- **Mini batch:** There are three steps.
 - Sample neighbourhood
 - Aggregate feature information from neighbours
 - Predict graph context and label using aggregated information



- Towards large-scale **GraphSAGE**:
 - Sampling mini-batch (sample target nodes as a mini-batch)
 - Sampling a fixed size set for each target nodes



- Any differentiable function that maps set of vectors to a single vector.



$$\mathbf{h}_v^k = \sigma \left(\left[\mathbf{A}_k \cdot \text{AGG}(\{\mathbf{h}_u^{k-1}, \forall u \in N(v)\}), \mathbf{B}_k \mathbf{h}_v^{k-1} \right] \right)$$

- How to aggregate information from neighbourhood

Concatenate neighbor embedding
and self embedding

$$\mathbf{h}_v^k = \sigma \left([\mathbf{W}_k \cdot \text{AGG}(\{\mathbf{h}_u^{k-1}, \forall u \in N(v)\}), \mathbf{B}_k \mathbf{h}_v^{k-1}] \right)$$

Generalized aggregation

- Mean

$$\text{AGG} = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|}$$

- Pooling

Element-wise mean/max

$$\text{AGG} = \gamma(\{\mathbf{Q}\mathbf{h}_u^{k-1}, \forall u \in N(v)\})$$

- LSTM

$$\text{AGG} = \text{LSTM}([\mathbf{h}_u^{k-1}, \forall u \in \pi(N(v))])$$

Algorithm 1: GraphSAGE embedding generation (i.e., forward propagation) algorithm

Input : Graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$; input features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; depth K ; weight matrices $\mathbf{W}^k, \forall k \in \{1, \dots, K\}$; non-linearity σ ; differentiable aggregator functions $\text{AGGREGATE}_k, \forall k \in \{1, \dots, K\}$; neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$

Output : Vector representations \mathbf{z}_v for all $v \in \mathcal{V}$

```

1  $\mathbf{h}_v^0 \leftarrow \mathbf{x}_v, \forall v \in \mathcal{V}$ ;
2 for  $k = 1 \dots K$  do
3   for  $v \in \mathcal{V}$  do
4      $\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$ ;
5      $\mathbf{h}_v^k \leftarrow \sigma(\mathbf{W}^k \cdot \text{CONCAT}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}^k))$ 
6   end
7    $\mathbf{h}_v^k \leftarrow \mathbf{h}_v^k / \|\mathbf{h}_v^k\|_2, \forall v \in \mathcal{V}$ 
8 end
9  $\mathbf{z}_v \leftarrow \mathbf{h}_v^K, \forall v \in \mathcal{V}$ 
    
```

Generalized Aggregators:

- Mean aggregator (GCN)
- Pooling aggregator
- LSTM aggregator

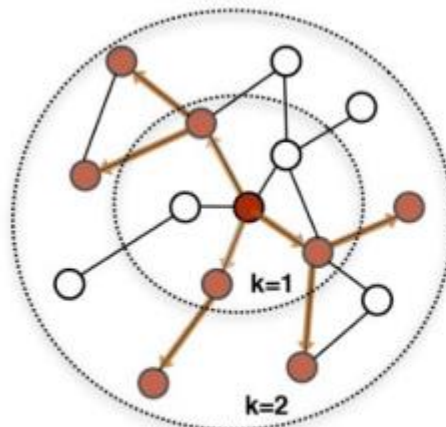
Use Concatenation instead of SUM

- Weighting factor in GraphSAGE
 - α_{uv} (importance) is defined explicitly based on the structure properties of graph.
 - All neighbors $u \in N(v)$ are equally important to node v

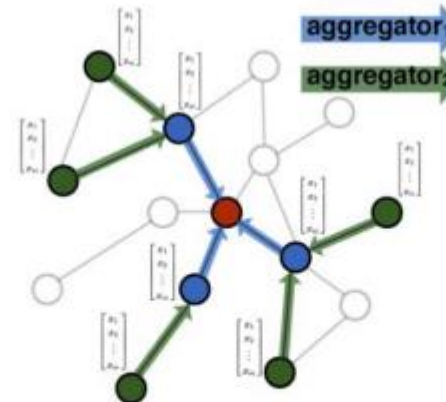
$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \boxed{\frac{\mathbf{h}_u^{k-1}}{|N(v)|}} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right)$$

Weighting
factor

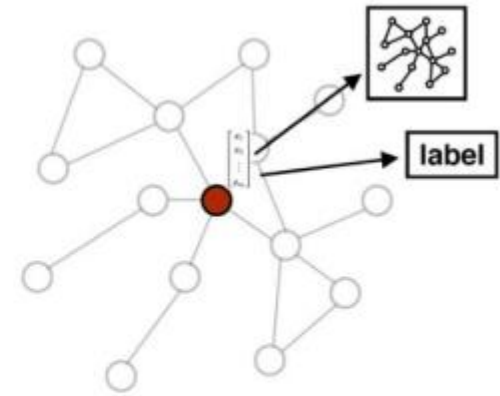
$$\alpha_{vu} = \frac{1}{|N(v)|}$$



1. Sample neighborhood



2. Aggregate feature information
from neighbors



3. Predict graph context and label
using aggregated information

- Pros:
 - Generalized aggregator.
 - Mini-batch training and fixed-size neighbor sampling.
- Cons:
 - Neighborhood expansion on deeper GNNs.
 - No guarantees for the sampling quality.

- DataLoader uses **NeighborSampler** to create mini-batches.
- Each mini-batch contains a node index and local graph information about that index. The key here is to sample local graph information for each mini-batch.
- **For example:**
 - Sampling neighboring nodes in each layer. sizes=[10, 5] means that 10 and 5 neighboring nodes are sampled in each layer.
 - We also use batch_size=32 to process 32 nodes in each batch.

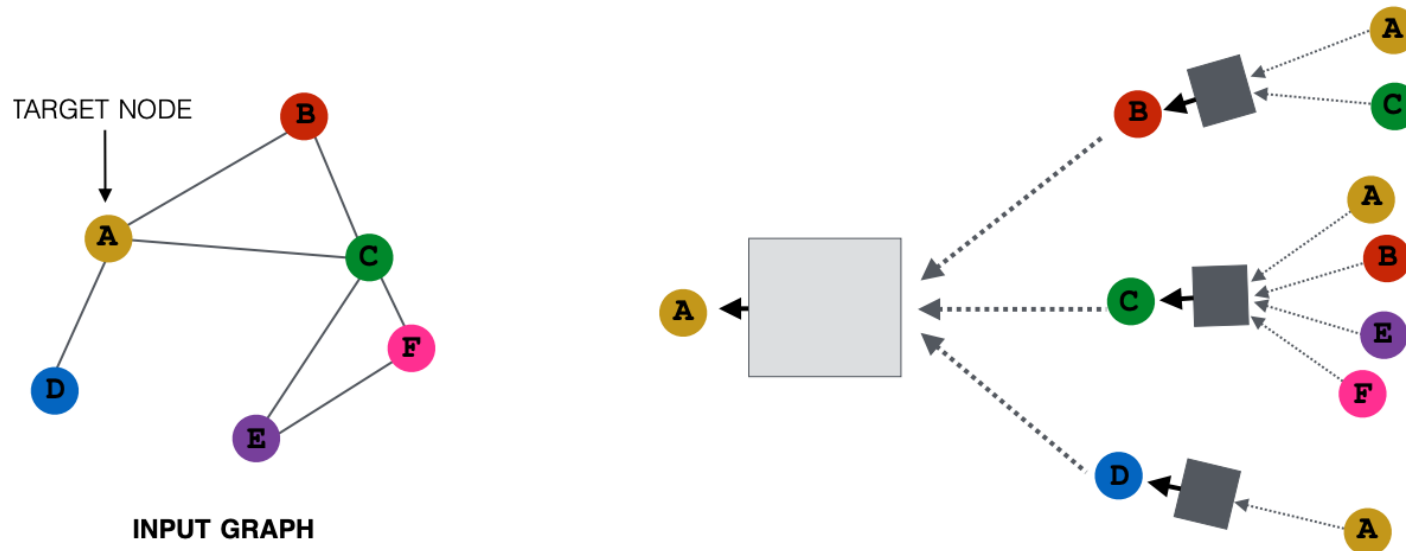
```
from torch_geometric.data import NeighborSampler

loader = NeighborSampler(data.edge_index, sizes=[10, 5], batch_size=32,
                        shuffle=True, num_nodes=data.num_nodes)
```

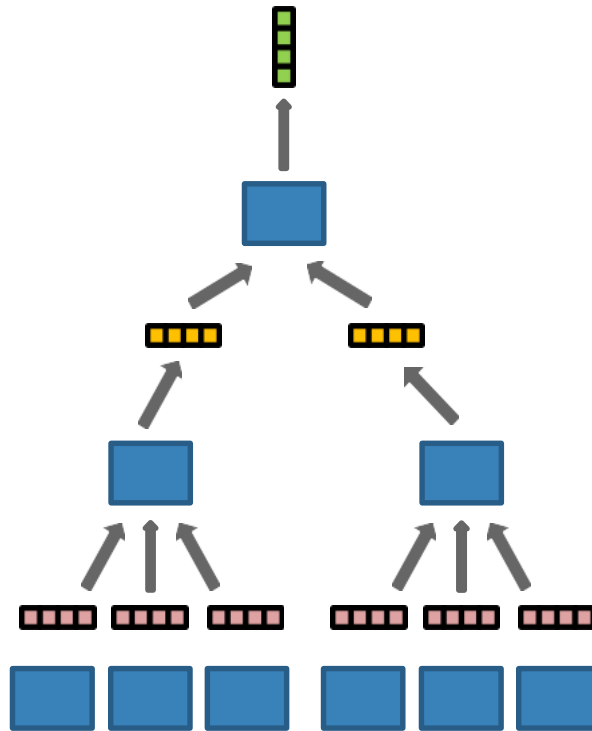
Lets do some example codes in the Sample code file

➤ **GCN is not trivial:**

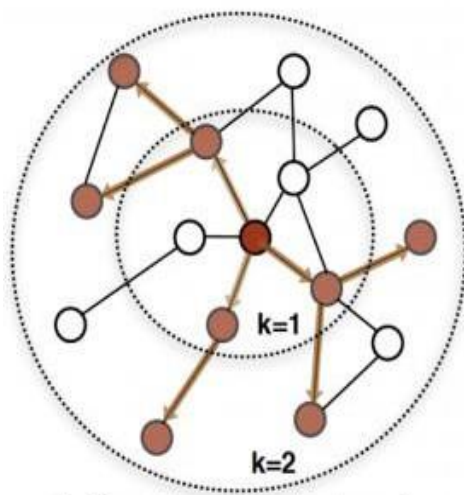
- In GCN, loss on a node not only depends on itself but all its neighbors
- This dependency brings difficulties when performing SGD on GCN



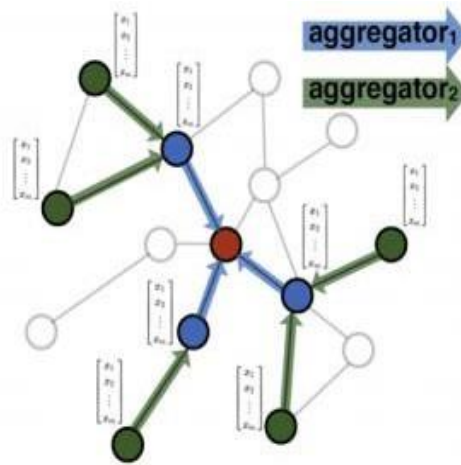
- Issues come from high computation costs
- Suppose we desire to calculate a target node's loss with a 2-layer GCN
- To obtain its final representation, needs all node embeddings in its 2-hop neighborhood
- **For example:** 9 nodes' embeddings needed but only get 1 loss (utilization: low)



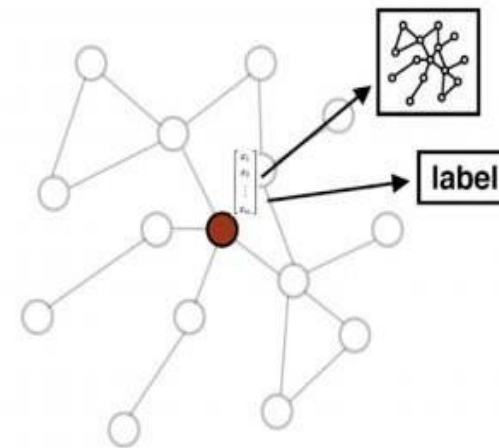
- **Idea:** subsample a smaller number of neighbors
 - For example, GraphSAGE (NeurIPS'17) considers a subset of neighbors per node
 - But it still suffers from recursive neighborhood expansion



1. Sample neighborhood



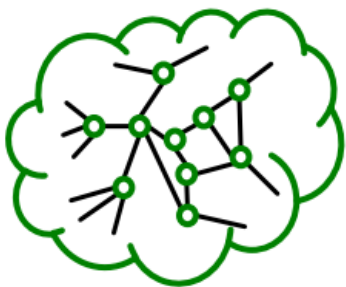
2. Aggregate feature information from neighbors



3. Predict graph context and label using aggregated information

- Idea: apply graph clustering algorithm (e.g., METIS) to identify dense subgraphs.

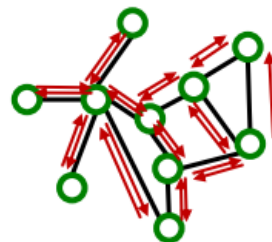
Large graph



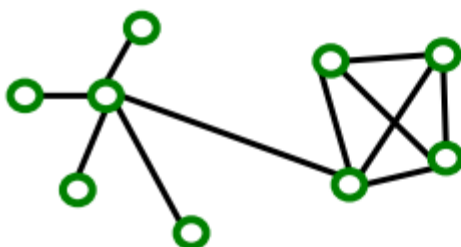
Sampled subgraph
(small enough to
be put on a GPU)



Layer-wise
node embeddings
update on the GPU



Original graph



Subgraphs (both 4-node induced subgraph)

Left



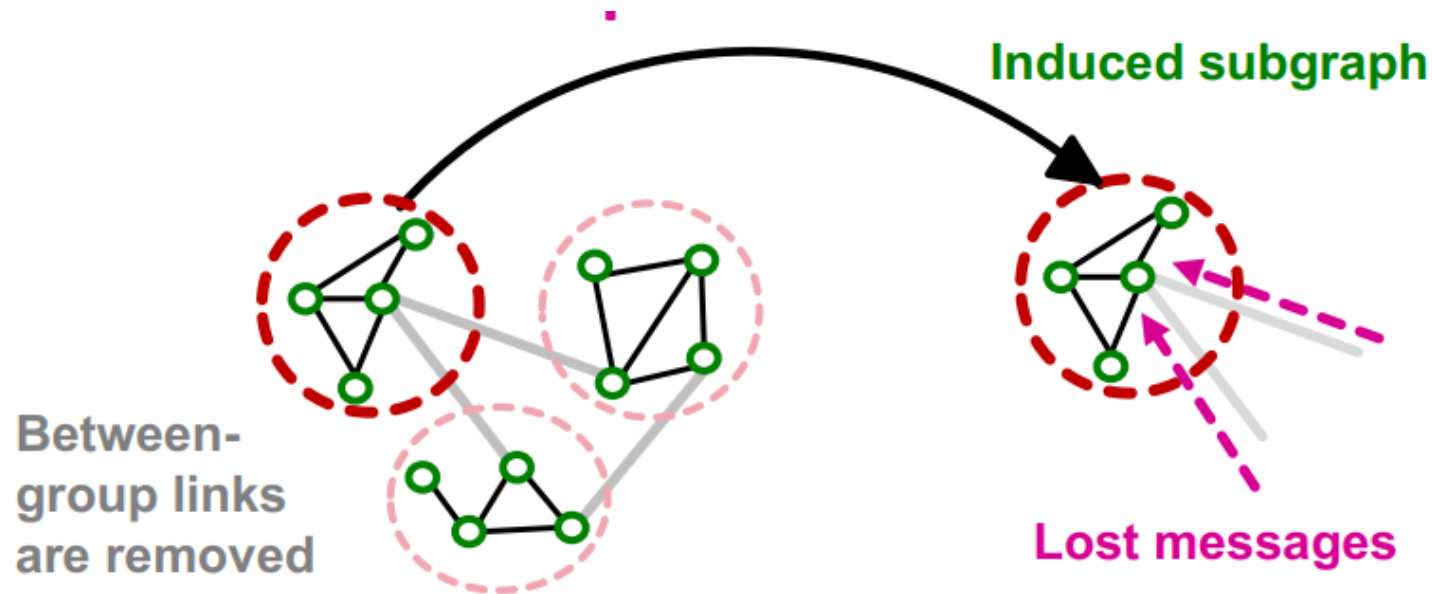
v.s.

Right



➤ **Problems:**

- The induced subgraph removes between group links.
- As a result, messages from other groups will be lost during message passing, which could hurt the GNN's performance.

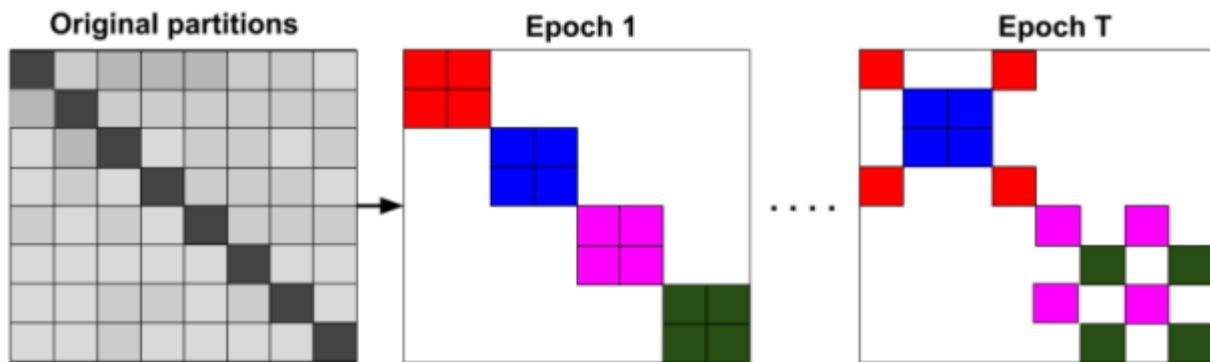
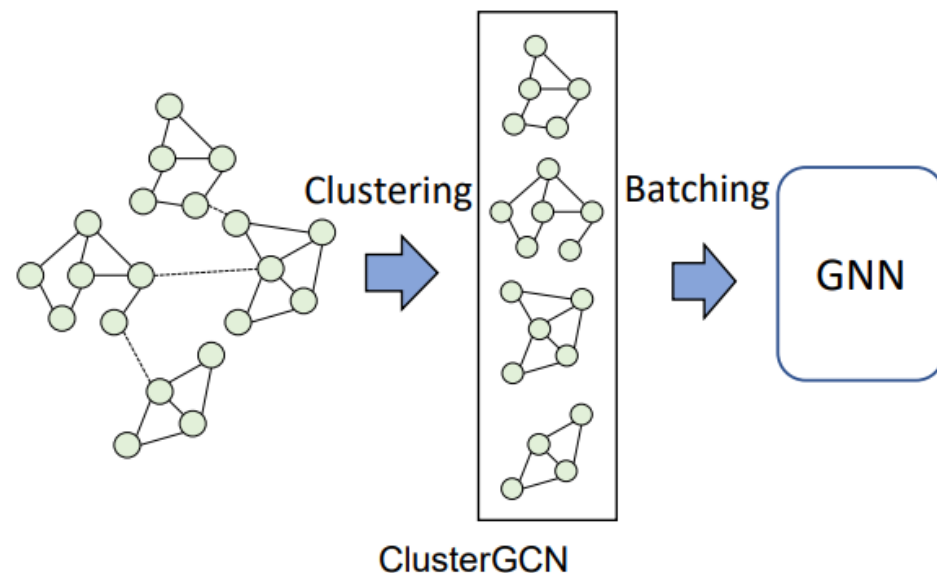


- Extract small clusters based efficient clustering algorithms.

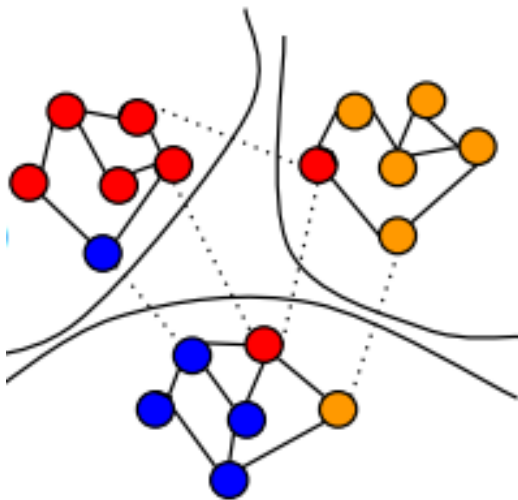
$$\bar{G} = [G_1, \dots, G_c] = [\{\mathcal{V}_1, \mathcal{E}_1\}, \dots, \{\mathcal{V}_c, \mathcal{E}_c\}],$$

$$\bar{A} = \begin{bmatrix} A_{11} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A_{cc} \end{bmatrix}, \Delta = \begin{bmatrix} 0 & \dots & A_{1c} \\ \vdots & \ddots & \vdots \\ A_{c1} & \dots & 0 \end{bmatrix},$$

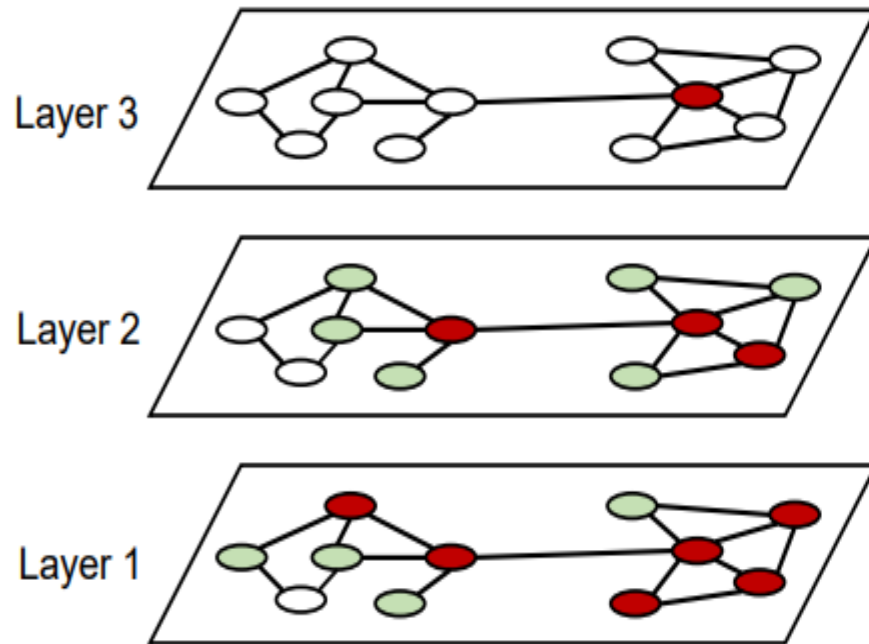
- Random batching at the subgraph level.



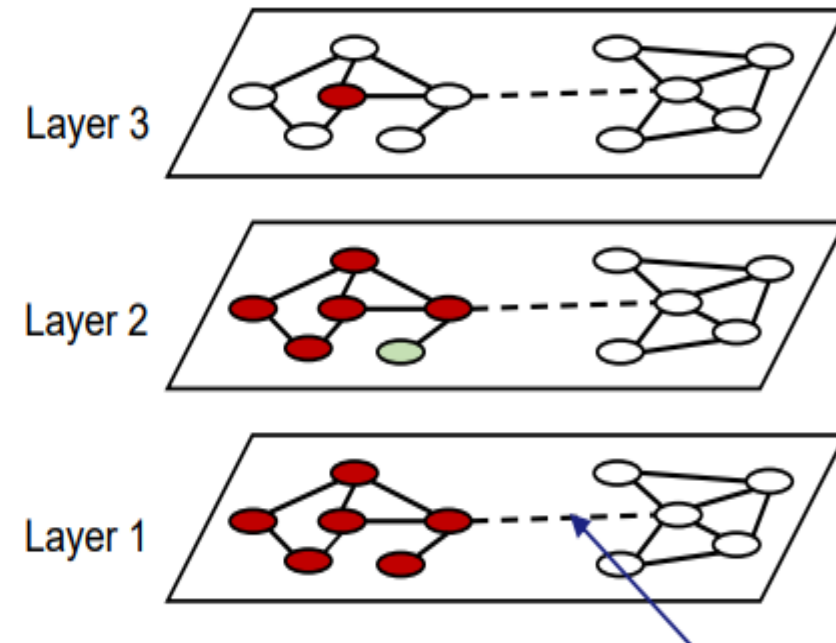
- **Idea:** apply graph clustering algorithm (e.g., METIS) to identify dense subgraphs.
- **Cluster-GCN**
 - Partition the graph into several clusters, remove between-cluster edges
 - Each subgraph is used as a mini-batch in SGD
 - Embedding utilization is optimal because nodes' neighbors stay within the cluster



- Neighbor expansion control:
Only consider the nodes in the same clusters



Fix-size neighbor sampling $S=2$

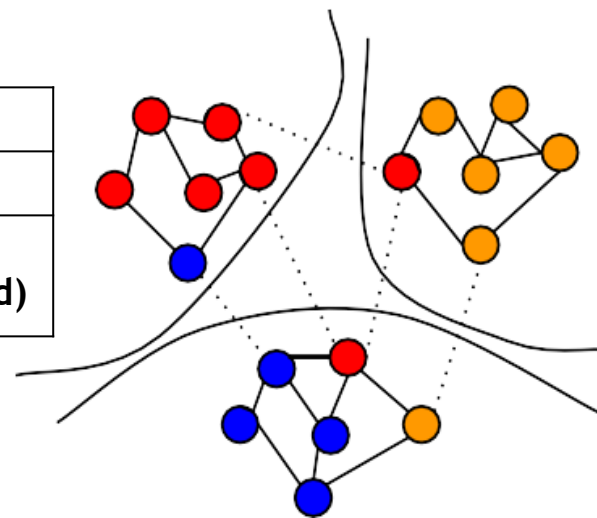


Only sample the nodes in the clusters

- Pros:
 - Good performance / Good memory usage.
 - Alleviate the neighborhood expansion problem in traditional mini-batch training.
- Cons:
 - Empirical results without analyzing the sampling quality.

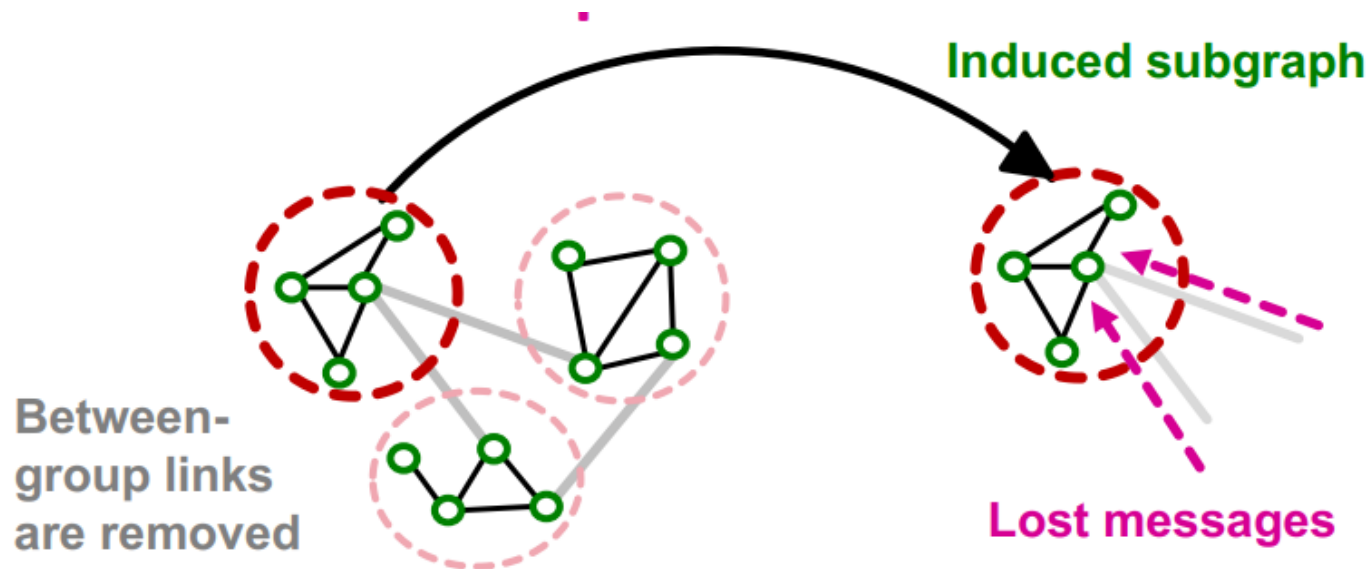
- Partition the graph into several clusters, remove between-cluster edges
- Each subgraph is used as a mini-batch in SGD
- Embedding utilization is optimal because nodes' neighbors stay within the cluster
- Even though 20% edges are removed, the accuracy of GCN model remains similar

| CiteSeer | Random partitioning | Graph partitioning |
|---------------------|---------------------|------------------------------|
| 1 (no partitioning) | 72.0 | 72.0 |
| 100 partitions | 46.1 | 71.5 (~20% edges removed) |

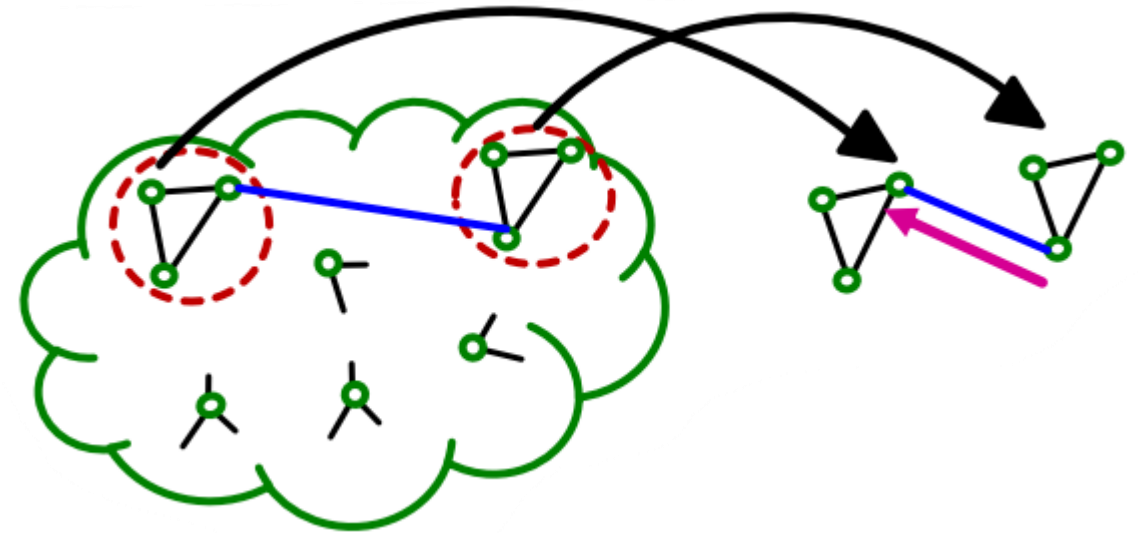
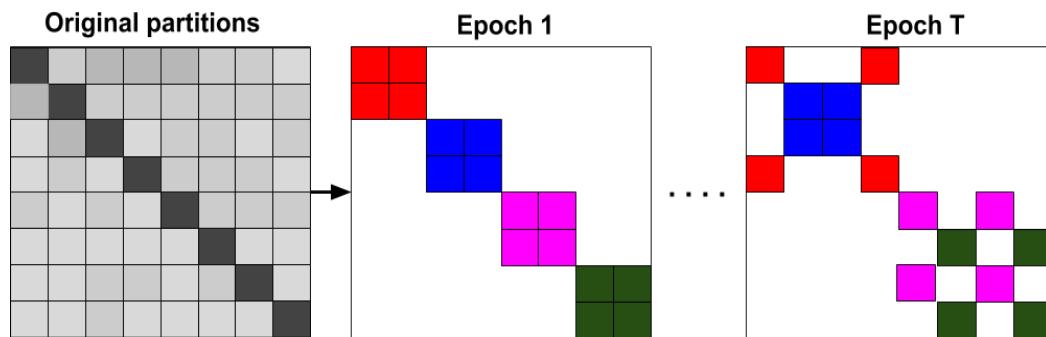


➤ Problems:

- The induced subgraph removes between group links.
- As a result, messages from other groups will be lost during message passing, which could hurt the GNN's performance.



- **Multiple Clusters**
 - A randomly select multiple clusters as a batch has been proposed.
- **Two advantages:**
 - Balance label distribution within a batch
 - Recover some missing edges between-cluster



Mini-batch training:

- For each mini-batch, **randomly sample a set of q node groups**: $\{V_{t_1}, \dots, V_{t_q}\} \subset \{V_1, \dots, V_C\}$.
- Aggregate all nodes across the sampled node groups**: $V_{aggr} = V_{t_1} \cup \dots \cup V_{t_q}$
- Extract the **induced subgraph** $G_{aggr} = (V_{aggr}, E_{aggr})$, where $E_{aggr} = \{(u, v) \mid u, v \in V_{aggr}\}$
 - E_{aggr} also includes between-group edges!

Algorithm 1: Cluster GCN

Input: Graph A , feature X , label Y ;

Output: Node representation \tilde{X}

- | | | | |
|--------------------------|---|---|--|
| Nodes Clustering | { | 1 | Partition graph nodes into c clusters $\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_c$ by METIS; |
| Random training clusters | { | 2 | for $iter = 1, \dots, max_iter$ do |
| | | 3 | Randomly choose q clusters, t_1, \dots, t_q from \mathcal{V} without replacement; |
| | | 4 | Form the subgraph \tilde{G} with nodes $\tilde{\mathcal{V}} = [\mathcal{V}_{t_1}, \mathcal{V}_{t_2}, \dots, \mathcal{V}_{t_q}]$ and links $A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}$; |
| Train and optimize | { | 5 | Compute $g \leftarrow \nabla \mathcal{L}_{A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}}$ (loss on the subgraph $A_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}}$); |
| | | 6 | Conduct Adam update using gradient estimator g |
| | | 7 | Output: $\{W_l\}_{l=1}^L$ |

- Clusters/partitions a graph data object into multiple subgraphs

```
from torch_geometric.datasets import Reddit
from torch_geometric.loader import ClusterData, ClusterLoader, NeighborLoader
from torch_geometric.nn import SAGEConv
```

```
dataset = Reddit('../data/Reddit')
data = dataset[0]
```

```
cluster_data = ClusterData(data, num_parts=1500, recursive=False,
                           save_dir=dataset.processed_dir)
train_loader = ClusterLoader(cluster_data, batch_size=20, shuffle=True,
                             num_workers=12)
```

```
subgraph_loader = NeighborLoader(data, num_neighbors=[-1], batch_size=1024,
                                 shuffle=False, num_workers=12)
```

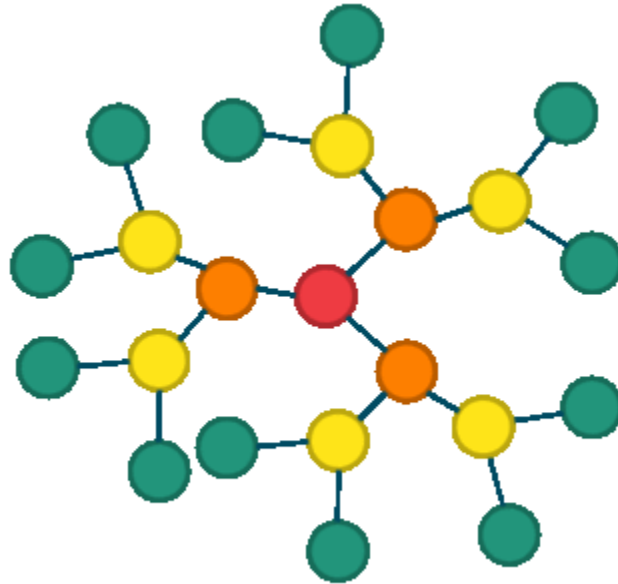
data: The graph data object
num_parts: The number of partitions.
Recursive: multilevel recursive bisection instead of multilevel k-way partitioning

ClusterLoader:
merges partitioned subgraphs and their between-cluster links from a large-scale graph data object to form a mini-batch.

NeighborLoader:
data loader that performs neighbour sampling. This loader allows for mini-batch training of GNNs on large-scale graphs where full-batch training is not feasible.

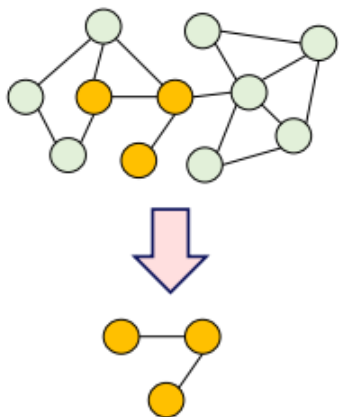
Lets do some example codes in the Sample code file

- "Neighbor explosion", GCNs rely on aggregating neighbor information to update nodes.
- It is very intuitive that the more layers of GCNs, the more neighbors need to be considered when updating nodes; while the number of layers is fixed
- The higher the average degree, the more neighbors need to be considered.



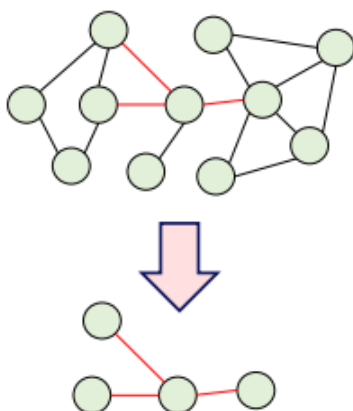
- Directly sample a subgraph for mini-batch training according to subgraph sampler.
- Sampler construction

Node sampler



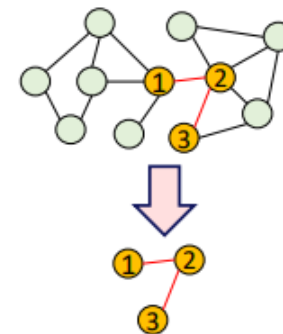
Uniformly sample nodes.

Edge sampler



Sample edge with probability $p_{u,v} \propto 1/d_u + 1/d_v$

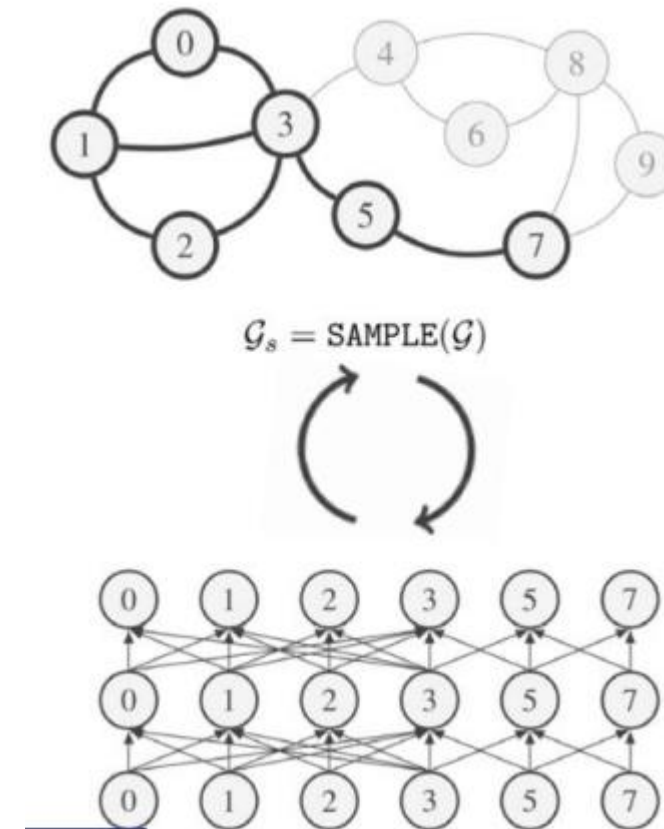
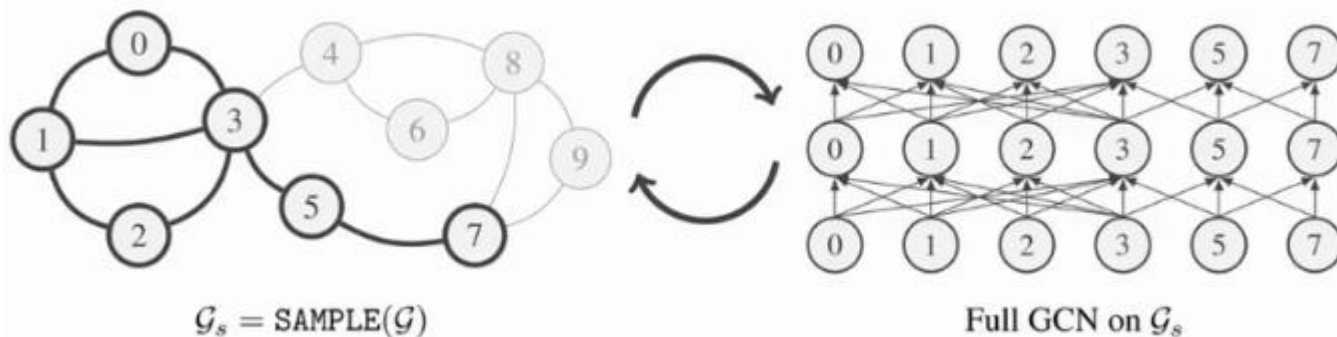
Random walk sampler



Sample edge with probability $p_{u,v} \propto B_{u,v} + B_{v,u}$

- $B_{u,v}$: the probability of a random walk to start at u and end at v in L hops.

- Sample a small subgraph, then build a complete GNN
- Constant neighborhood size
- Not an I.I.D data sampler
- **Why?**
 - Popular users will be sampled more frequently, i.e., influencers, ...
- **How?**
 - Normalize aggreg and mini batch loss by edge/node sampling frequency.



- How to eliminate the bias introduced by the sampler?

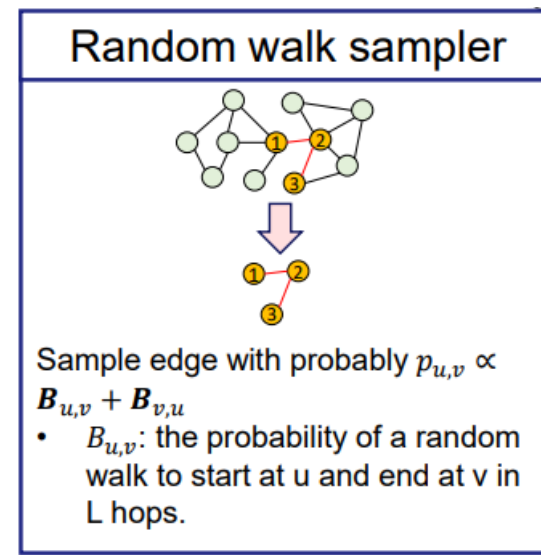
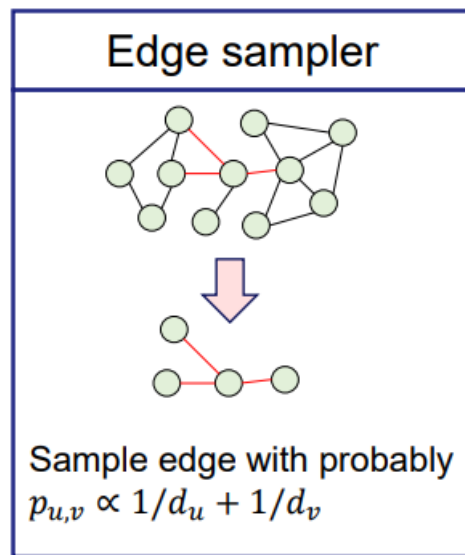
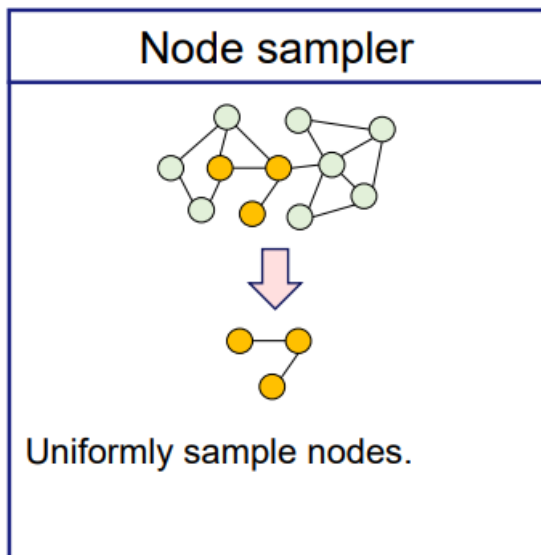
- Loss normalization:

$$\mathcal{L}_{\text{batch}} = \sum_{v \in G_s} L_v / \lambda_v, \lambda_v = |V| p_v.$$

- Aggregation normalization:

$$a(u, v) = p_{u,v} / p_v$$

p_v : the probability of a node $v \in V$ being sampled.
 p_{uv} : the probability of an edge $u, v \in E$ being sampled.



```
from torch_geometric.datasets import Flickr
from torch_geometric.loader import GraphSAINTRandomWalkSampler
from torch_geometric.nn import GraphConv
from torch_geometric.typing import WITH_TORCH_SPARSE
from torch_geometric.utils import degree

loader = GraphSAINTRandomWalkSampler(data, batch_size=6000, walk_length=2,
                                     num_steps=5, sample_coverage=100,
                                     save_dir=dataset.processed_dir,
                                     num_workers=4)
```

Given a graph, this class samples nodes and constructs subgraphs that can be processed in a mini – batch fashion.

Lets do some example codes in the Sample code file



네트워크 과학연구실
NETWORK SCIENCE LAB



가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

